



# Expression of Endogenous Genes by Non-homologous Recombination of a Vector Construct With Cellular DNA

## *Background of the Invention*

### 5 *Field of the Invention*

The field of the invention is activating gene expression or causing over-expression of a gene by recombination methods *in situ*. The invention relates to expressing an endogenous gene in a cell at levels higher than those normally found in the cell. Expression of the gene is activated or increased following integration, by non-homologous or illegitimate recombination, of a regulatory sequence that activates expression of the gene. The method allows the identification and expression of genes undiscoverable by current methods since no target sequence is necessary for integration. Thus, gene products associated with human disease and development are obtainable from genes that have not been sequenced and indeed, whose existence is unknown, as well as from well-characterized genes. The methods provide gene products from such genes for therapeutic and diagnostic purposes.

### *Related Art*

Identification and over-expression of novel genes associated with human disease is an important step towards developing new therapeutic drugs. Current approaches to creating libraries of cells for protein over-expression are based on the production and cloning of cDNA. Thus, in order to identify a new gene using this approach, the gene must be expressed in the cells that were used to make the library. The gene also must be expressed at sufficient levels to be adequately represented in the library. This is problematic because many genes are expressed only in very low quantities, in a rare population of cells, or during short developmental periods.

Furthermore, because of the large size of some mRNAs, it is difficult or impossible to produce full length cDNA molecules capable of expressing the biologically active protein. Lack of full-length cDNA molecules has also been observed for small mRNAs and is thought to be related to sequences in the message that are difficult to produce by reverse transcription or that are unstable during propagation in bacteria. As a result, even the most complete cDNA libraries express only a fraction of the entire set of possible genes.

Finally, many cDNA libraries are produced in bacterial vectors. Use of these vectors to express biologically active mammalian proteins is severely limited since most mammalian proteins do not fold correctly and/or are improperly glycosylated in bacteria.

Therefore, a method for creating a more representative library for protein expression, capable of facilitating faithful expression of biologically active proteins, would be extremely valuable.

Current methods for over-expressing proteins involve cloning the gene of interest and placing it, in a construct, next to a suitable promoter/enhancer, polyadenylation signal, and splice site, and introducing the construct into an appropriate host cell.

An alternative approach involves the use of homologous recombination to activate gene expression by targeting a strong promoter or other regulatory sequence to a previously identified gene.

WO 90/14092 describes *in situ* modification of genes, in mammalian cells, encoding proteins of interest. This application describes single-stranded oligonucleotides for site-directed modification of genes encoding proteins of interest. A marker may also be included. However, the methods are limited to providing an oligonucleotide sequence substantially homologous to a target site. Thus, the method requires knowledge of the site required for activation by site-directed modification and homologous recombination. Novel genes are not discoverable by such methods.

WO 91/06667 describes methods for expressing a mammalian gene *in situ*. With this method, an amplifiable gene is introduced next to a target gene by homologous recombination. When the cell is then grown in the appropriate medium, both the amplifiable gene and the target gene are amplified and there is enhanced expression of the target gene. As above, methods of introducing the amplifiable gene are limited to homologous recombination, and are not useful for activating novel genes whose sequence (or existence) is unknown.

WO 91/01140 describes the inactivation of endogenous genes by modification of cells by homologous recombination. By these methods, homologous recombination is used to modify and inactivate genes and to produce cells which can serve as donors in gene therapy.

WO 92/20808 describes methods for modifying genomic target sites *in situ*. The modifications are described as being small, for example, changing single bases in DNA. The method relies upon genomic modification using homologous DNA for targeting.

WO 92/19255 describes a method for enhancing the expression of a target gene, achieved by homologous recombination in which a DNA sequence is integrated into the genome or large genomic fragment. This modified sequence can then be transferred to a secondary host for expression. An amplifiable gene can be integrated next to the target gene so that the target region can be amplified for enhanced expression. Homologous recombination is necessary to this targeted approach.

WO 93/09222 describes methods of making proteins by activating an endogenous gene encoding a desired product. A regulatory region is targeted by homologous recombination and replacing or disabling the region normally associated with the gene whose expression is desired. This disabling or replacement causes the gene to be expressed at levels higher than normal.

WO 94/12650 describes a method for activating expression of and amplifying an endogenous gene *in situ* in a cell, which gene is not expressed or is not expressed at desired levels in the cell. The cell is transfected with

exogenous DNA sequences which repair, alter, delete, or replace a sequence present in the cell or which are regulatory sequences not normally functionally linked to the endogenous gene in the cell. In order to do this, DNA sequences homologous to genomic DNA sequences at a preselected site are used to target the endogenous gene. In addition, amplifiable DNA encoding a selectable marker can be included. By culturing the homologously recombinant cells under conditions that select for amplification, both the endogenous gene and the amplifiable marker are co-amplified and expression of the gene increased.

WO 95/31560 describes DNA constructs for homologous recombination. The constructs include a targeting sequence, a regulatory sequence, an exon, and an unpaired splice donor site. The targeting is achieved by homologous recombination of the construct with genomic sequences in the cell and allows the production of a protein *in vitro* or *in vivo*.

WO 96/29411 describes methods using an exogenous regulatory sequence, an exogenous exon, either coding or non-coding, and a splice donor site introduced into a preselected site in the genome by homologous recombination. In this application, the introduced DNA is positioned so that the transcripts under control of the exogenous regulatory region include both the exogenous exon and endogenous exons present in either the thrombopoietin, DNase I, or  $\beta$ -interferon genes, resulting in transcripts in which the exogenous and endogenous exons are operably linked. The novel transcription units are produced by homologous recombination.

U.S. Patent No. 5,272,071 describes the transcriptional activation of transcriptionally silent genes in a cell by inserting a DNA regulatory element capable of promoting the expression of a gene normally expressed in that cell. The regulatory element is inserted so that it is operably linked to the normally silent gene. The insertion is accomplished by means of homologous recombination by creating a DNA construct with a segment of the normally silent gene (the target DNA) and the DNA regulatory element used to induce the desired transcription.

U.S. Patent No. 5, 578,461 discusses activating expression of mammalian target genes by homologous recombination. A DNA sequence is integrated into the genome or a large genomic fragment to enhance the expression of the target gene. The modified construct can then be transferred to a secondary host. An amplifiable gene can be integrated adjacent to the target gene so that the target region is amplified for enhanced expression.

Both of the above approaches (construction of an over-expressing construct by cloning or by homologous recombination *in vivo*) require the gene to be cloned and sequenced before it can be over-expressed. Furthermore, using homologous recombination, the genomic sequence and structure must also be known.

Unfortunately, many genes have not yet been identified and/or sequenced. Thus, a method for over-expressing a gene of interest, whether or not it has been previously cloned, and whether or not its sequence and structure are known, would be useful.

### *Summary of the Invention*

The invention is, therefore, generally directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell. The method does not require previous knowledge of the sequence of the endogenous gene or even of the existence of the gene.

The cell containing the vector is screened for expression of the gene.

The cell over-expressing the gene can be cultured *in vitro* so as to produce desired amounts of the gene product of the endogenous gene that has been activated or whose expression has been increased. The gene product can then be isolated and purified to use, for example, in protein therapy or drug discovery.

Alternatively, the cell expressing the desired gene product can be allowed to express the gene product *in vivo*.

The vector construct can consist essentially of the transcriptional regulatory sequence.

5           The invention is also directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence and an amplifiable marker into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

10           The cell containing the vector is screened for over-expression of the gene.

          The cell over-expressing the gene is cultured such that amplification of the endogenous gene is obtained. The cell can then be cultured *in vitro* so as to produce desired amounts of the gene product of the amplified endogenous gene that has been activated or whose expression has been increased. The gene  
15           product can then be isolated and purified.

          Alternatively, following amplification, the cell can be allowed to express the endogenous gene and produce desired amounts of the gene product *in vivo*.

          The vector construct can consist essentially of the transcriptional regulatory sequence and the amplifiable marker.

20           It is to be understood, however, that any vector used in the methods described herein can include an amplifiable marker. Thereby, amplification of both the vector and the DNA of interest (i.e., containing the over-expressed gene) occurs in the cell, and further enhanced expression of the endogenous gene is obtained. Accordingly, methods can include a step in which the endogenous  
25           gene is amplified.

          The invention is also directed to methods for over-expressing an endogenous gene in a cell comprising introducing a vector containing a transcriptional regulatory sequence and an unpaired splice donor sequence into the cell, allowing the vector to integrate into the genome of the cell by

non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

The cell containing the vector is screened for expression of the gene.

5 The cell over-expressing the gene can be cultured *in vitro* so as to produce desirable amounts of the gene product of the endogenous gene whose expression has been activated or increased. The gene product can then be isolated and purified.

Alternatively, the cell can be allowed to express the desired gene product *in vivo*.

10 The vector construct can consist essentially of the transcriptional regulatory sequence and the splice donor sequence.

Any of the vector constructs used in the methods described herein can also include a secretion signal sequence. The secretion signal sequence is arranged in the construct so that it will be operably linked to the activated endogenous protein. Thereby, secretion of the protein of interest occurs in the cell, and purification of that protein is facilitated. Accordingly, methods can include a step in which the protein expression product is secreted from the cell.

15 The invention also encompasses cells made by any of the above methods. The invention encompasses cells containing the vector constructs, cells in which the vector constructs have integrated, and cells which are over-expressing desired gene products from an endogenous gene, over-expression being driven by the introduced transcriptional regulatory sequence.

The cells can be isolated and cloned.

25 The methods can be carried out in any cell of eukaryotic origin, such as fungal, plant or animal. Preferred embodiments include vertebrates and particularly mammals, and more particularly, humans.

A single cell made by the methods described above can over-express a single gene or more than one gene. More than one gene in a cell can be activated by the integration of a single type of construct into multiple locations in the genome. Similarly, more than one gene in a cell can be activated by the

30

integration of multiple constructs (i.e., more than one type of construct) into multiple locations in the genome. Therefore, a cell can contain only one type of vector construct or different types of constructs, each capable of activating an endogenous gene.

5           The invention is also directed to methods for making the cells described above by one or more of the following: introducing one or more of the vector constructs; allowing the introduced construct(s) to integrate into the genome of the cell by non-homologous recombination; allowing over-expression of one or more endogenous genes in the cell; and isolating and cloning the cell.

10           The invention also encompasses methods for using the cells described above to over-express a gene that has been characterized (for example, sequenced), uncharacterized (for example, a gene whose function is known but which has not been cloned or sequenced), or a gene whose existence was, prior to over-expression, unknown. The cells can be used to provide desired amounts  
15 of a gene product *in vitro* or *in vivo*. The gene product can then be isolated and purified if desired. It can be purified by cell lysis or from the growth medium (as when the vector contains a secretion signal sequence).

          The invention also encompasses libraries of cells made by the above described methods. A library can encompass all of the clones from a single  
20 transfection experiment or a subset of clones from a single transfection experiment. The subset can over-express the same gene or more than one gene, for example, a class of genes. The transfection can have been done with a single construct or with more than one construct.

          A library can also be formed by combining all of the recombinant cells  
25 from two or more transfection experiments, by combining one or more subsets of cells from a single transfection experiment or by combining subsets of cells from separate transfection experiments. The resulting library can express the same gene, or more than one gene, for example, a class of genes. Again, in each of these individual transfections, a unique construct or more than one construct can  
30 be used.



Libraries can be formed from the same cell type or different cell types.

The invention is also directed to methods for making libraries by selecting various subsets of cells from the same or different transfection experiments.

5 The invention accordingly is also directed to methods of using libraries of cells to over-express endogenous genes. The library is screened for the expression of the gene and cells are selected that express the desired gene product. The cell can then be used to purify the gene product for subsequent use. Expression in the cell can occur by culturing the cell *in vitro* or by allowing the cell to express the gene *in vivo*.

10 In preferred embodiments of the invention, the methods include a process wherein the expression product is purified. In highly preferred embodiments, the cells expressing the endogenous gene product are cultured so as to produce amounts of gene product feasible for commercial application, and especially diagnostic and therapeutic and drug discovery uses.

15 Any of the methods can further comprise introducing double-strand breaks into the genomic DNA in the cell prior to or simultaneously with vector integration.

20 The invention also encompasses novel vector constructs for activating gene expression or over-expressing a gene through non-homologous recombination. The novel construct lacks homologous targeting sequences. That is, it does not contain nucleotide sequences that target host cell DNA and promote homologous recombination at the target site, causing over-expressing of a cellular gene via the introduced transcriptional regulatory sequence.

25 Novel vector constructs include a vector containing a transcriptional regulatory sequence operably linked to an unpaired splice donor sequence and further contains an amplifiable marker.

30 Novel vector constructs include constructs with a transcriptional regulatory sequence operably linked to a translational start codon, a signal secretion sequence, and an unpaired splice donor site; constructs with a transcriptional regulatory sequence, operably linked to a translation start codon,

an epitope tag, and an unpaired splice donor site; constructs containing a transcriptional regulatory sequence operably linked to a translational start codon, a signal sequence and an epitope tag, and an unpaired splice donor site; constructs containing a transcriptional regulatory sequence operably linked to a translation  
5 start codon, a signal secretion sequence, an epitope tag, and a sequence-specific protease site, and an unpaired splice donor site.

The vector construct can contain a selectable marker for recombinant host cell selection. Alternatively, selection can be effected by phenotypic selection for a trait provided by the activated endogenous gene product.

10 These vectors, and indeed any of the vectors disclosed herein, and obvious variants recognized by one of ordinary skill in the art, can be used in any of the methods described herein to form any of the compositions producible by these methods.

The transcriptional regulatory sequence includes, but is not limited to, a  
15 promoter. In preferred embodiments, the promoter is a viral promoter. In highly preferred embodiments, the viral promoter is the cytomegalovirus immediate early promoter. In alternative embodiments, the promoter is a cellular, non-viral promoter or inducible promoter.

The transcriptional regulatory sequence also includes, but is not limited  
20 to, an enhancer. In preferred embodiments, the enhancer is a viral enhancer. In highly preferred embodiments, the viral enhancer is the cytomegalovirus immediate early enhancer. In alternative embodiments, the enhancer is a cellular non-viral enhancer.

In preferred embodiments of the methods described herein, the vector  
25 construct is or contains linear RNA or DNA.

### ***Brief Description of the Figures***

FIG. 1. Schematic diagram of gene activation events described herein. The activation construct is transfected into cells and allowed to integrate into the

host cell chromosomes at DNA breaks. If breakage occurs upstream of a gene of interest (e.g., EPO), and the appropriate activation construct integrates at the break such that its regulatory sequence becomes operably linked to the gene of interest, activation of the gene will occur. Transcription and splicing produce a chimeric RNA molecule containing exonic sequences from the activation construct and from the endogenous gene. Subsequent translation will result in the production of the protein of interest. Following isolation of the ecombinant cell, gene expression can be further enhanced via gene amplification.

**FIG. 2.** Schematic diagram of non-translated activation constructs. The arrows denote promoter sequences. The exonic sequences are shown as open boxes and the splice donor sequence is indicated by S/D. Construct numbers corresponding to the description below are shown on the left. The selectable and amplifiable markers are not shown.

**FIG. 3.** Schematic diagram of translated activation constructs. The arrows denote promoter sequences. The exonic sequences are shown as open boxes and the splice donor sequence is indicated by S/D. The translated, signal peptide, epitope tag, and protease cleavage sequences are shown in the legend below the constructs. Construct numbers corresponding to the description below are shown on the left. The selectable and amplifiable markers are not shown.

**FIG. 4.** Schematic diagram of an activation construct capable of activating endogenous genes.

### *Detailed Description of the Preferred Embodiments*

There are great advantages of gene activation by non-homologous recombination over other gene activation procedures. Unlike previous methods of protein over-expression, the methods described herein do not require that the gene of interest be cloned (isolated from the cell). Nor do they require any knowledge of the DNA sequence or structure of the gene to be over-expressed (i.e., the sequence of the ORF, introns, exons, or upstream and downstream

regulatory elements) or knowledge of a gene's expression patterns (i.e., tissue specificity, developmental regulation, etc.). Furthermore, the methods do not require any knowledge pertaining to the genomic organization of the gene of interest (i.e., the intron and exon structure).

5           The methods of the present invention thus involve vector constructs that do not contain target nucleotide sequences for homologous recombination. A target sequence allows homologous recombination of vector DNA with cellular DNA at a predetermined site on the cellular DNA, the site having homology for sequences in the vector, the homologous recombination at the predetermined site  
10           resulting in the introduction of the transcriptional regulatory sequence into the genome and the subsequent endogenous gene activation.

          The method of the present invention does not involve integration of the vector at predetermined sites.

          The vectors described herein do not contain target sequences. A target  
15           sequence is a sequence on the vector that has homology with a sequence or sequences within the gene to be activated or upstream of the gene to be activated, the upstream region being up to and including the first functional splice acceptor site on the same coding strand of the gene of interest, and by means of which  
20           homology the transcriptional regulatory sequence that activates the gene of interest is integrated into the genome of the cell containing the gene to be activated. In the case of an enhancer integration vector for activating an endogenous gene, the vector does not contain homology to any sequence in the genome upstream or downstream of the gene of interest (or within the gene of  
25           interest) for a distance extending as far as enhancer function is operative.

          The methods, therefore, are capable of identifying new genes that have  
30           been or can be missed using conventional and currently available cloning techniques. By using the constructs and methodology described herein, unknown and/or uncharacterized genes can be rapidly identified and over-expressed to produce proteins. The proteins have use as, among other things, human therapeutics and diagnostics and targets for drug discovery.

The methods are also capable of producing over-expression of known and/or characterized genes for *in vitro* or *in vivo* protein production.

5 A "known" gene relates to the level of characterization of a gene. The invention allows expression of genes that have been characterized as well as of genes that have not been characterized. Different levels of characterization are possible. These include detailed characterization, such as cloning, DNA, RNA, and/or protein sequencing, and relating the regulation and function of the gene to the cloned sequence (e.g., recognition of promoter and enhancer sequences, functions of the open reading frames, introns, and the like). Characterization can be less detailed, such as having mapped a gene and related function, or having a partial amino acid or nucleotide sequence, or having purified a protein and ascertained a function. Characterization may be minimal, as when a nucleotide or amino acid sequence is known or a protein has been isolated but the function is unknown. Alternatively, a function may be known but the associated protein or nucleotide sequence is not known or is known but is not related to the function. Finally, there may be no characterization in that both the existence of the gene and its function are not known. The invention allows expression of any gene at any of these or other specific degrees of characterization.

20 Many different proteins can be activated or over-expressed by a single activation construct and in a single set of transfections. Thus, a single cell or different cells in a set of transfectants (library) can over-express more than one protein following transfection with the same or different constructs. Previous activation methods require a unique construct to be created for each gene to be activated.

25 Further, many different integration sites adjacent to a single gene can be created and tested simultaneously using a single construct. This allows rapid determination of the optimal genomic location of the activation construct for protein expression.

30 Using previous methods, the 5' end of the gene of interest had to be extensively characterized with respect to sequence and structure. For each

activation construct to be produced, an appropriate targeting sequence had to be isolated. Usually, this must be an isogenic sequence isolated from the same person or laboratory strain of animal as the cells to be activated. In some cases, this DNA may be 50 kb or more from the gene of interest. Thus, production of each targeting construct required an arduous amount of cloning and sequencing of the endogenous gene. Since sequence and structure information is not required for the methods of the invention, unknown genes and genes with uncharacterized upstream regions can be activated.

This is made possible using methods of *in situ* gene activation using non-homologous recombination of exogenous DNA sequences with cellular DNA.

DNA molecules can recombine to redistribute their genetic content by several different and distinct mechanisms, including homologous recombination, site-specific recombination, and non-homologous/illegitimate recombination. Homologous recombination involves recombination between stretches of DNA that are highly similar in sequence. It has been demonstrated that homologous recombination involves pairing between the homologous sequences along their length prior to redistribution of the genetic material. The exact site of crossover can be at any point in the homologous segments. The efficiency of recombination is proportional to the length of homologous targeting sequence (Hopé, *Development* 113:399 (1991); Reddy *et al.*, *J. Virol.* 65:1507 (1991)), the degree of sequence identity between the two recombining sequences (von Melchner *et al.*, *Genes Dev.* 6:919 (1992)), and the ratio of homologous to non-homologous DNA present on the construct (Letson, *Genetics* 117:759 (1987)).

Site-specific recombination, on the other hand, involves the exchange of genetic material at a predetermined site, designated by specific DNA sequences. In this reaction, a protein recombinase binds to the recombination signal sequences, creates a strand scission, and facilitates DNA strand exchange. *Cre/Lox* recombination is an example of site specific recombination.

Non-homologous/illegitimate recombination involves the joining (exchange or redistribution) of genetic material that does not share significant

sequence homology and does not occur at site-specific recombination sequences. Examples of non-homologous recombination include integration of exogenous DNA into chromosomes at non-homologous sites, chromosomal translocations and deletions, DNA end-joining, double strand break repair of chromosome ends, bridge-breakage fusion, and concatemerization of transfected sequences. In most cases, non-homologous recombination is thought to occur through the joining of "free DNA ends." Free ends are DNA molecules that contain an end capable of being joined to a second DNA end either directly, or following repair or processing. The DNA end may consist of a 5' overhang, 3' overhang, or blunt end.

As used herein, retroviral insertion and other transposition reactions are loosely considered forms of non-homologous recombination. These reactions do not involve the use of homology between the recombining molecules. Furthermore, unlike site-specific recombination, these types of recombination reactions do not occur between discrete sites. Instead, a specific protein/DNA complex is required on only one of the recombination partners (i.e., the retrovirus or transposon), with the second DNA partner (i.e., the cellular genome) usually being relatively non-specific. As a result, these "vectors" do not integrate into the cellular genome in a targeted fashion, and therefore they can be used to deliver the activation construct according to the present invention.

Vector constructs useful for the methods described herein contain a transcriptional regulatory sequence that undergoes non-homologous recombination with genomic sequences in a cell to over-express an endogenous gene in that cell. The vector construct lacks homologous targeting sequences. That is, it does not contain DNA sequences that target host cell DNA and promote homologous recombination at the target site, causing over-expressing of a cellular gene via the introduced transcriptional regulatory sequence.

The invention is generally directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence into the cell, allowing the vector to integrate

into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell. The method does not require previous knowledge of the sequence of the endogenous gene or even of the existence of the gene. Where the sequence of the gene to be activated is known,  
5 however, the constructs can be engineered to contain the proper configuration of vector elements (e.g., location of the start codon, addition of codons present in the first exon of the the endogenous gene, and the proper reading frame) to achieve maximal overexpression and/or the appropriate protein sequence.

The cell containing the vector is screened for expression of the gene.

10 The cell over-expressing the gene can be cultured *in vitro* so as to produce desired amounts of the gene product of the endogenous gene that has been activated or whose expression has been increased. The gene product can then be isolated and purified to use, for example, in protein therapy or drug discovery.

Alternatively, the cell expressing the desired gene product can be allowed  
15 to express the gene product *in vivo*.

The vector construct can consist essentially of the transcriptional regulatory sequence.

Alternatively, the vector construct can consist essentially of the transcriptional regulatory sequence and the amplifiable marker.

20 The invention, therefore, is also directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence and an amplifiable marker into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

25 The cell containing the vector is screened for over-expression of the gene.

The cell over-expressing the gene is cultured such that amplification of the endogenous gene is obtained. The cell can then be cultured *in vitro* so as to produce desired amounts of the gene product of the amplified endogenous gene that has been activated or whose expression has been increased. The gene  
30 product can then be isolated and purified.



Alternatively, following amplification, the cell can be allowed to express the endogenous gene and produce desired amounts of the gene product *in vivo*.

The vector construct can consist essentially of the transcriptional regulatory sequence and the splice donor sequence.

5 The invention, therefore, is also directed to methods for over-expressing an endogenous gene in a cell comprising introducing a vector containing a transcriptional regulatory sequence and an unpaired splice donor sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous  
10 gene in the cell.

The cell containing the vector is screened for expression of the gene.

The cell over-expressing the gene can be cultured *in vitro* so as to produce desirable amounts of the gene product of the endogenous gene whose expression has been activated or increased. The gene product can then be isolated and  
15 purified.

Alternatively, the cell can be allowed to express the desired gene product *in vivo*.

The vector construct can consist essentially of a transcriptional regulatory sequence operably linked to an unpaired splice donor sequence and also  
20 containing an amplifiable marker.

Other activation vectors include constructs with a transcriptional regulatory sequence and an exonic sequence containing a start codon; a transcriptional regulatory sequence and an exonic sequence containing a translational start codon and a secretion signal sequence; constructs with a  
25 transcriptional regulatory sequence and an exonic sequence containing a translation start codon, and an epitope tag; constructs containing a transcriptional regulatory sequence and an exonic sequence containing a translational start codon, a signal sequence and an epitope tag; constructs containing a transcriptional regulatory sequence and an exonic sequence with a translation start  
30 codon, a signal secretion sequence, an epitope tag, and a sequence-specific

protease site. In each of the above constructs, the exon on the construct is located immediately upstream of an unpaired splice donor site.

5 The constructs can also contain a regulatory sequence, a selectable marker lacking a poly A signal, an internal ribosome entry site (ires), and an unpaired splice donor site (FIG. 4). A start codon, signal secretion sequence, epitope tag, and/or a protease cleavage site may optionally be included between the ires and the unpaired splice donor sequence. When this construct integrates upstream of a gene, the selectable marker will be efficiently expressed since a poly A site will be supplied by the endogenous gene. In addition the downstream gene will also  
10 be expressed since the ires will allow protein translation to initiate at the downstream open reading frame (i.e. the endogenous gene). Thus, the message produced by this activation construct will be polycistronic. The advantage of this construct is that integration events that do not occur near genes and in the proper orientation will not produce a drug resistant colony. The reason for this is that  
15 without a poly A tail (supplied by the endogenous gene), the neomycin resistance gene will not express efficiently. By reducing the number of nonproductive integration events, the complexity of the library can be reduced without affecting its coverage (the number of genes activated), and this will facilitate the screening process.

20 In another embodiment of this construct, *cre-lox* recombination sequences can be included between the regulatory sequence and the *neo* start codon and between the ires and the unpaired splice donor site (between the ires and the start codon, if present). Following isolation of cells that have activated the gene of interest, the *neo* gene and ires can be removed by transfecting the cells with a  
25 plasmid encoding the *cre* recombinase. This would eliminate the production of the polycistronic message and allow the endogenous gene to be expressed directly from the regulatory sequence on the integrated activation construct. Use of *Cre* recombination to facilitate deletion of genetic elements from mammalian chromosomes has been described (Gu *et al.*, *Science* 265:103 (1994); Sauer, *Meth. Enzymology* 225:890-900 (1993)).  
30

Thus, constructs useful in the methods described herein include, but are not limited to, the following (See also Figures 1-4):

- 1) Construct with a regulatory sequence and an exon lacking a translation start codon.
- 5      2) Construct with a regulatory sequence and an exon lacking a translation start codon followed by a splice donor site.
- 3) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 10     4) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 5) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 15     6) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 7) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 20     8) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 25     9) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 30     10) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.

- 11) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 5 12) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 10 13) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 15 14) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 20 15) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 25 16) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 30 17) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 18) Construct with a regulatory sequence linked to a selectable marker, followed by an internal ribosome entry site, and an unpaired splice donor site.

- 19) Construct 18 in which a cre/lox recombination signal is located between  
a) the regulatory sequence and the open reading frame of the selectable  
marker and b) between the ires and the unpaired splice donor site.
- 20) Construct with a regulatory sequence operably linked to an exon  
containing green fluorescent protein lacking a stop codon, followed by an  
unpaired splice donor site.

5

It is to be understood, however, that any vector used in the methods  
described herein can include an amplifiable marker. Accordingly, methods can  
include a step in which the endogenous gene is amplified. Placement of an  
10 amplifiable marker on the activation construct results in the juxtaposition of the  
gene of interest and the amplifiable marker in the activated cell. Once the  
activated cell has been isolated, expression can be further increased by selecting  
for cells containing an increased copy number of the locus containing both the  
gene of interest and the activation construct. This can be accomplished by  
15 selection methods known in the art.

Examples of amplifiable markers include dihydrofolate reductase,  
adenosine deaminase, aspartate transcarbamylase, dihydro-orotase, and carbamyl  
phosphate synthase.

20

It is also understood that any of the constructs described herein may  
contain a eukaryotic viral origin of replication, either in place of, or in  
conjunction with an amplifiable marker. The presence of the viral origin of  
replication allows the integrated vector and adjacent endogenous gene to be  
isolated as an episome and/or amplified to high copy number upon introduction  
of the appropriate viral replication protein. Examples of useful viral origins  
25 include, but are not limited to, SV40 ori and EBV ori P.

The invention encompasses embodiments in which the constructs  
disclosed herein consist essentially of the components specifically described for  
these constructs.

It is also understood that the above constructs are examples of constructs useful in the methods described herein, but that the invention encompasses functional equivalents of such constructs.

5 The term "vector" is understood to generally refer to the vehicle by which the nucleotide sequence is introduced into the cell. It is not intended to be limited to any specific sequence. The vector could itself be the nucleotide sequence that activates the endogenous gene or could contain the sequence that activates the endogenous gene. Thus, the vector could be simply a linear or circular polynucleotide containing essentially only those sequences necessary for  
10 activation, or could be these sequences in a larger polynucleotide or other construct such as a DNA or RNA viral genome, a whole virion, or other biological construct used to introduce the critical nucleotide sequences into a cell.

The vector can contain DNA sequences that exist in nature or that have been created by genetic engineering or synthetic processes.

15 The construct, upon nonhomologous integration into the genome of a cell, can activate expression of an endogenous gene. Expression of the endogenous gene may result in production of full length protein, or in production of a truncated biologically active form of the endogenous protein, depending on the integration site (e.g., upstream region versus intron 2). The activated gene may  
20 be a known gene (e.g., previously cloned or characterized) or unknown gene (previously not cloned or characterized). The function of the gene may be known or unknown.

Examples of proteins with known activities include, but are not limited to, cytokines, growth factors, enzymes, structural proteins, cell surface receptors,  
25 intracellular receptors, hormones, antibodies, and transcription factors. Specific examples of known proteins that can be produced by this method include, but are not limited to, erythropoietin, insulin, growth hormone, glucocerebrosidase, tissue plasminogen activator, granulocyte-colony stimulating factor, granulocyte/macrophage colony stimulating factor, interferon  $\alpha$ , interferon  $\beta$ ,  
30 interferon  $\gamma$ , interleukin-2, interleukin-6, interleukin-11, interleukin-12, TGF  $\beta$ ,

blood clotting factor V, blood clotting factor VII, blood clotting factor VIII, blood clotting factor IX, blood clotting factor X, TSH- $\beta$ , bone growth factor 2, bone growth factor-7, tumor necrosis factor, alpha-1 antitrypsin, anti-thrombin III, leukemia inhibitory factor, glucagon, Protein C, protein kinase C, macrophage colony stimulating factor, stem cell factor, follicle stimulating hormone  $\beta$ , urokinase, nerve growth factors, insulin-like growth factors, insulinotropin, parathyroid hormone, lactoferrin, complement inhibitors, platelet derived growth factor, keratinocyte growth factor, neurotrophin-3, thrombopoietin, chorionic gonadotropin, thrombomodulin, alpha glucosidase, epidermal growth factor, FGF, macrophage-colony stimulating factor, interleukin 4, interleukin 10, and cell surface receptors for each of the above growth factors, hormones, and cytokines.

One of the advantages of the method described herein is that virtually any gene can be activated. However, since genes have different genomic structures, including different intron/exon boundaries and locations of start codons, a variety of activation constructs is provided to activate the maximum number of different genes within a population of cells.

These constructs can be transfected separately into cells to produce libraries. Each library contains cells with a unique set of activated genes. Some genes will be activated by several different activation constructs. In addition, portions of a gene can be activated to produce truncated, biologically active proteins. Truncated proteins can be produced, for example, by integration of an activation construct into introns or exons in the middle of an endogenous gene rather than upstream of the second exon.

Use of different constructs also allows the activated gene to be modified to contain new sequences. For example, a secretion signal sequence can be included on the activation construct to facilitate the secretion of the activated gene. In some cases, depending on the intron/exon structure or the gene of interest, the secretion signal sequence can replace all or part of the signal sequence of the endogenous gene. In other cases, the signal sequence will allow a protein which is normally located intracellularly to be secreted.

The regulatory sequence on the vector can be a constitutive promoter. Alternatively, the promoter may be inducible. Use of inducible promoters will allow low basal levels of activated protein to be produced by the cell during routine culturing and expansion. The cells may then be induced to produce large amounts of the desired proteins, for example, during manufacturing or screening. Examples of inducible promoters include, but are not limited to, the tetracycline inducible promoter and the metallothionein promoter.

The regulatory sequence on the vector can be a tissue specific promoter.

The regulatory sequence on the vector can be an enhancer.

The regulatory sequence on the vector can be isolated from cellular or viral genomes. Examples of cellular regulatory sequences include, but are not limited to, regulatory elements from the actin gene, metallothionein I gene, immunoglobulin genes, casein I gene, serum albumin gene, collagen gene, globin genes, laminin gene, spectrin gene, ankyrin gene, sodium/potassium ATPase gene, and tubulin gene. Examples of viral regulatory sequences include, but are not limited to, regulatory elements from CMV immediate early gene, adenovirus late genes, SV40 genes, retroviral LTRs, and Herpesvirus genes. Typically, regulatory sequences contain binding sites for transcription factors such as NF-kB, SP-1, TATA binding protein, AP-1, and CAAT binding protein. Functionally, the regulatory sequence is defined by its ability to promote, enhance, or otherwise alter transcription of an endogenous gene.

In preferred embodiments, the regulatory sequence is a viral promoter. In highly preferred embodiments, the promoter is the cytomegalovirus immediate early gene promoter. In alternative embodiments, the regulatory element is a cellular, non-viral promoter.

In preferred embodiments, the regulatory element contains an enhancer. In highly preferred embodiments, the enhancer is the cytomegalovirus immediate early gene enhancer. In alternative embodiments, the enhancer is a cellular, non-viral enhancer.



The transcriptional regulatory sequence can be comprised of scaffold-attachment regions or matrix attachment sites, negative regulatory elements, and transcription factor binding sites. Regulatory sequences can also include locus control regions.

5           The invention encompasses the use of retrovirus transcriptional regulatory sequences, e.g., long terminal repeats. Where these are used, however, they are not necessarily linked to any retrovirus sequence that materially affects the function of the transcriptional regulatory sequence as a promoter or enhancer of transcription of the endogenous gene to be activated (i.e., the cellular gene with  
10           which the transcriptional regulatory sequence recombines to activate).

          The construct may contain a regulatory sequence which is not operably linked to exonic sequences on the vector. For example, when the regulatory element is an enhancer, it can integrate near an endogenous gene (e.g., upstream, downstream, or in an intron) and stimulate expression of the gene from its  
15           endogenous promoter. By this mechanism of activation, exonic sequences from the vector are absent in the transcript of the activated gene.

          Alternatively, the regulatory element may be operably linked to an exon. The exon may be a naturally occurring sequence or may be non-naturally occurring (e.g., produced synthetically). To activate endogenous genes lacking  
20           a start codon in their first exon (e.g., follicle stimulating hormone- $\beta$ ), a start codon is preferably omitted from the exon on the vector. To activate endogenous genes containing a start codon in the first exon (e.g., erythropoietin and growth hormone), the exon on the vector preferably contains a start codon, usually ATG and preferably an efficient translation initiation site (Kozak, *J. Mol Biol.* 196: 947  
25           (1987)). The exon may contain additional codons following the start codon. These codons may be derived from a naturally occurring gene or may be non-naturally occurring (e.g., synthetic). The codons may be the same as the codons present in the first exon of the endogenous gene to be activated. Alternatively, the codons may be different than the codons present in the first  
30           exon of the endogenous gene. For example, the codons may encode an epitope

tag, signal secretion sequence, transmembrane domain, selectable marker, or screenable marker. Optionally, an unpaired splice donor site may be present immediately 3' of the exonic sequence. When the structure of the gene to be activated is known, the splice donor site should be placed adjacent to the vector exon in a location such that the codons in the vector will be in frame with the codons of the second exon of the endogenous gene following splicing. When the structure of the endogenous gene to be activated is not known, separate constructs, each containing a different reading frame, are used.

Operably linked is defined as a configuration that allows transcription through the designated sequence(s). For example, a regulatory sequence that is operably linked to an exonic sequence indicates that the exonic sequence is transcribed. When a start codon is present on the vector, operably linked also indicates that the open reading frame from the vector exon is in frame with the open reading frame of the endogenous gene. Following nonhomologous integration, the regulatory sequence (e.g., a promoter) on the vector becomes operably linked to an endogenous gene and facilitates transcription initiation, at a site generally referred to as a CAP site. Transcription proceeds through the exonic elements on the vector (and, if present, through the start codon, open reading frame, and/or unpaired splice donor site), and through the endogenous gene. The primary transcript produced by this operable linkage is spliced to create a chimeric transcript containing exonic sequences from both the vector and the endogenous gene. This transcript is capable of producing the endogenous protein when translated.

An exon or "exonic sequence" is defined as any transcribed sequence that is present in the mature RNA molecule. The exon on the vector may contain untranslated sequences, for example, a 5' untranslated region. Alternatively, or in conjunction with the untranslated sequences, the exon may contain coding sequences such as a start codon and open reading frame. The open reading frame can encode naturally occurring amino acid sequences or non-naturally occurring amino acid sequences (e.g., synthetic codons). The open reading frame may also

encode a signal secretion sequence, epitope tag, exon, selectable marker, screenable marker, or nucleotides that function to allow the open reading frame to be preserved when spliced to an endogenous gene.

5 Splicing of primary transcripts, the process by which introns are removed, is directed by a splice donor site and a splice acceptor site, located at the 5' and 3' ends of introns, respectively. The consensus sequence for splice donor sites is (A/C)AG GURAGU (where R represents a purine nucleotide) with nucleotides in positions 1-3 located in the exon and nucleotides GURAGU located in the intron.

10 An unpaired splice donor site is defined herein as a splice donor site present on the activation construct without a downstream splice acceptor site. When the vector is integrated by nonhomologous recombination into a host cell's genome, the unpaired splice donor site becomes paired with a splice acceptor site from an endogenous gene. The splice donor site from the vector, in conjunction  
15 with the splice acceptor site from the endogenous gene, will then direct the excision of all of the sequences between the vector splice donor site and the endogenous splice acceptor site. Excision of these intervening sequences removes sequences that interfere with translation of the endogenous protein.

20 The terms upstream and downstream, as used herein, are intended to mean in the 5' or in the 3' direction, respectively, relative to the coding strand. The term "upstream region" of a gene is defined as the nucleic acid sequence 5' of its second exon (relative to the coding strand) up to and including the last exon of the first adjacent gene having the same coding strand. Functionally, the upstream region is any site 5' of the second exon of an endogenous gene capable of  
25 allowing a nonhomologously integrated vector to become operably linked to the endogenous gene.

The vector construct can contain a selectable marker to facilitate the identification and isolation of cells containing a nonhomologously integrated activation construct. Examples of selectable markers include genes encoding  
30 neomycin resistance (neo), hypoxanthine phosphoribosyl transferase (HPRT),

puromycin (pac), dihydro-orotase glutamine synthetase (GS), histidine D (his D), carbamyl phosphate synthase (CAD), dihydrofolate reductase (DHFR), multidrug resistance 1 (mdr1), aspartate transcarbamylase, xanthine-guanine phosphoribosyl transferase (gpt), and adenosine deaminase (ada).

5           Alternatively, the vector can contain a screenable marker, in place of or in addition to, the selectable marker. A screenable marker allows the cells containing the vector to be isolated without placing them under drug or other selective pressures. Examples of screenable markers include genes encoding cell surface proteins, fluorescent proteins, and enzymes. The vector containing cells  
10           may be isolated, for example, by FACS using fluorescently-tagged antibodies to the cell surface protein or substrates that can be converted to fluorescent products by a vector encoded enzyme.

          Alternatively, selection can be effected by phenotypic selection for a trait provided by the endogenous gene product. The activation construct, therefore,  
15           can lack a selectable marker other than the "marker" provided by the endogenous gene itself. In this embodiment, activated cells can be selected based on a phenotype conferred by the activated gene. Examples of selectable phenotypes include cellular proliferation, growth factor independent growth, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell,  
20           muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), expression of cell surface receptors/proteins, gain or loss of cell-cell adhesion, migration, and cellular activation (e.g., resting versus activated T cells).

          A selectable marker may also be omitted from the construct when  
25           transfected cells are screened for gene activation products without selecting for the stable integrants. This is particularly useful when the efficiency of stable integration is high.

          The vector may contain an amplifiable marker to allow for selection of cells containing increased copies of the integrated vector and the adjacent  
30           activated endogenous gene. Examples of amplifiable markers include

dihydrofolate reductase (DHFR), adenosine deaminase (ada), dihydro-orotase glutamine synthetase (GS), and carbamyl phosphate synthase (CAD).

5 The vector may contain eukaryotic viral origins of replication useful for gene amplification. These origins may be present in place of, or in conjunction with, an amplifiable marker.

The vector may also contain genetic elements useful for the propagation of the construct in micro-organisms. Examples of useful genetic elements include microbial origins of replication and antibiotic resistance markers.

10 These vectors, and any of the vectors disclosed herein, and obvious variants recognized by one of ordinary skill in the art, can be used in any of the methods described above to form any of the compositions producible by those methods.

15 Nonhomologous integration of the construct into the genome of a cell results in the operable linkage between the regulatory elements from the vector and the exons from an endogenous gene. In preferred embodiments, the insertion of the vector regulatory sequences is used to upregulate expression of the endogenous gene. Upregulation of gene expression includes converting a transcriptionally silent gene to a transcriptionally active gene. It also includes enhancement of gene expression for genes that are already transcriptionally  
20 active, but produce protein at levels lower than desired. In other embodiments, expression of the endogenous gene may be affected in other ways such as downregulation of expression, creation of an inducible phenotype, or changing the tissue specificity of expression.

25 Cells produced by this method can be used to produce protein *in vitro* (e.g., for use as a protein therapeutic) or *in vivo* (e.g., for use in cell therapy).

The invention also encompasses cells made by any of the above methods. The invention encompasses cells containing the vector constructs, cells in which the vector constructs have integrated, and cells which are over-expressing desired gene products from an endogenous gene. over-expression being driven by the  
30 introduced transcriptional regulatory sequence.

Cells used in this invention can be derived from any eukaryotic species and can be primary, secondary, or immortalized. Furthermore, the cells can be derived from any tissue in the organism. Examples of useful tissues from which cells can be isolated and activated include, but are not limited to, liver, kidney, spleen, bone marrow, thymus, heart, muscle, lung, brain, testes, ovary, islet, intestinal, bone marrow, skin, bone, gall bladder, prostate, bladder, embryos, and the immune and hematopoietic systems. Cell types include fibroblast, epithelial, neuronal, stem, and follicular. However, any cell or cell type can be used to activate gene expression using this invention.

The methods can be carried out in any cell of eukaryotic origin, such as fungal, plant or animal. Preferred embodiments include vertebrates and particularly mammals, and more particularly, humans.

The construct can be integrated into primary, secondary, or immortalized cells. Primary cells are cells that have been isolated from a vertebrate and have not been passaged. Secondary cells are primary cells that have been passaged, but are not immortalized. Immortalized cells are cell lines that can be passaged, apparently indefinitely.

In preferred embodiments, the cells are immortalized cell lines. Examples of immortalized cell lines include, but are not limited to, HT1080, HeLa, Jurkat, 293 cells, KB carcinoma, T84 colonic epithelial cell line, Raji, Hep G2 hepatoma cell line, A2058 melanoma, U937 lymphoma, and WI38 fibroblast cell line, somatic cell hybrids, and hybridomas.

Cells used in this invention can be derived from any eukaryotic species. To overexpress endogenous human proteins, human cells are used. Similarly, to over-express endogenous bovine proteins, for example bovine growth hormone, bovine cells are used.

The cells can be derived from any tissue in the vertebrate organism. Examples of useful tissues from which cells can be isolated and activated include, but are not limited to, liver, kidney, spleen, bone marrow, thymus, heart, muscle, lung, brain, immune system, testes, ovary, islet, intestinal, bone marrow, skin,

bone, gall bladder, prostate, bladder, embryos, and hematopoietic. Cell types include fibroblast, epithelial, neuronal, stem, and follicular. However, any eukaryotic cell or cell type can be used to activate gene expression using this invention.

5           Any of the cells produced by any of the methods described are useful for screening for expression of a desired gene product and for providing desired amounts of a gene product that is over-expressed in the cell. The cells can be isolated and cloned.

10           Commercial growth and production conditions often vary from the conditions used to grow and prepare cells for analytical use (e.g., cloning, protein or nucleic acid sequencing, raising antibodies, X-ray crystallography analysis, enzymatic analysis, and the like). Scale up of cells for growth in roller bottles involves increase in the surface area on which cells can attach. Microcarrier beads are, therefore, often added to increase the surface area for commercial  
15           growth. Scale up of cells in spinner culture may involve large increases in volume. Five liters or greater can be required for both microcarrier and spinner growth. Depending on the inherent potency (specific activity) of the protein of interest, the volume can be as low as 1-10 liters. 10-15 liters is more common. However, up to 50-100 liters may be necessary and volume can be as high as  
20           10,000-15,000 liters. In some cases, higher volumes may be required. Cells can also be grown in large numbers of T flasks, for example 50-100.

          Despite growth conditions, protein purification on a commercial scale can also vary considerably from purification for analytic purposes. Protein purification in a commercial practical context can be initially the equivalent of  
25           10 liters of cells at approximately  $10^4$  cells/ml. Cell mass equivalent to begin protein purification can also be as high as 10 liters of cells at up to  $10^6$  or  $10^7$  cells/ml.

          Another commercial growth condition, especially when the ultimate product is used clinically, is cell growth in serum-free medium, by which is  
30           intended medium containing no serum or not in amounts that are required for cell

growth. This obviously avoids the undesired co-purification of toxic contaminants (e.g., viruses) or other types of contaminants, for example, proteins that would complicate purification.

5 A single cell made by the methods described above can over-express a single gene or more than one gene. More than one gene can be activated by the integration of a single construct or by the integration of multiple constructs in the same cell (i.e., more than one type of construct). Therefore, a cell can contain only one type of vector construct or different types of constructs, each capable of activating an endogenous gene.

10 The invention is also directed to methods for making the cells described above by one or more of the following: introducing one or more of the vector constructs; allowing the introduced construct(s) to integrate into the genome of the cell by non-homologous recombination; allowing over-expression of one or more endogenous genes in the cell; and isolating and cloning the cell.

15 The term "transfection" has been used herein for convenience when discussing introducing a polynucleotide into a cell. However, it is to be understood that the specific use of this term has been applied to generally refer to the *introduction* of the polynucleotide into a cell and is also intended to refer to the introduction by other methods described herein such as electroporation, liposome-mediated introduction, retrovirus-mediated introduction, and the like  
20 (as well as according to its own specific meaning).

The vector can be introduced into the cell by a number of methods known in the art. These include, but are not limited to, electroporation, calcium phosphate precipitation, DEAE dextran, lipofection, and receptor mediated  
25 endocytosis, polybrene, particle bombardment, and microinjection. Alternatively, the vector can be delivered to the cell as a viral particle (either replication competent or deficient). Examples of viruses useful for the delivery of nucleic acid include, but are not limited to, adenovirus, adeno-associated virus, retrovirus, Herpesviruses, and vaccinia virus.



Following transfection, the cells are cultured under conditions, as known in the art, suitable for nonhomologous integration between the vector and the host cell's genome. Cells containing the nonhomologously integrated vector can be further cultured under conditions, as known in the art, allowing expression of activated endogenous genes.

The vector construct can be introduced into cells on a single DNA construct or on separate constructs and allowed to concatemerize.

Whereas in preferred embodiments, the vector construct is a double-stranded DNA vector construct, vector constructs also include single-stranded DNA, combinations of single- and double-stranded DNA, single-stranded RNA, double-stranded RNA, and combinations of single- and double-stranded RNA. Thus, for example, the vector construct could be single-stranded RNA which is converted to cDNA by reverse transcriptase, the cDNA converted to double-stranded DNA, and the double-stranded DNA ultimately recombining with the host cell genome.

In preferred embodiments, the constructs are linearized prior to introduction into the cell. Linearization of the activation construct creates free DNA ends capable of reacting with chromosomal ends during the integration process. In general, the construct is linearized downstream of the regulatory element (and exonic and splice donor sequences, if present). Linearization can be facilitated by, for example, placing a unique restriction site downstream of the regulatory sequences and treating the construct with the corresponding restriction enzyme prior to transfection. While not required, it is advantageous to place a "spacer" sequence between the linearization site and the proximal most functional element (e.g., the unpaired splice donor site) on the construct. When present, the spacer sequence protects the important functional elements on the vector from exonucleolytic degradation during the transfection process. The spacer can be composed of any nucleotide sequence that does not change the essential functions of the vector as described herein.

Circular constructs can also be used to activate endogenous gene expression. It is known in the art that circular plasmids, upon transfection into cells, can integrate into the host cell genome. Presumably, DNA breaks occur in the circular plasmid during the transfection process, thereby generating free DNA ends capable of joining to chromosome ends. Some of these breaks in the construct will occur in a location that does not destroy essential vector functions (e.g., the break will occur downstream of the regulatory sequence), and therefore, will allow the construct to be integrated into a chromosome in a configuration capable of activating an endogenous gene. As described above, spacer sequences may be placed on the construct (e.g., downstream of the regulatory sequences). During transfection, breaks that occur in the spacer region will create free ends at a site in the construct suitable for activation of an endogenous gene following integration into the host cell genome.

The invention also encompasses libraries of cells made by the above described methods. A library can encompass all of the clones from a single transfection experiment or a subset of clones from a single transfection experiment. The subset can over-express the same gene or more than one gene, for example, a class of genes. The transfection can have been done with a single type of construct or with more than one type of construct.

A library can also be formed by combining all of the recombinant cells from two or more transfection experiments, by combining one or more subsets of cells from a single transfection experiment or by combining subsets of cells from separate transfection experiments. The resulting library can express the same gene, or more than one gene, for example, a class of genes. Again, in each of these individual transfections, a unique construct or more than one construct can be used.

Libraries can be formed from the same cell type or different cell types.

The library can be composed of a single type of cell containing a single type of activation construct which has been integrated into chromosomes at spontaneous DNA breaks or at breaks generated by radiation, restriction enzymes,

and/or DNA breaking agents, applied either together (to the same cells) or separately (applied to individual groups of cells and then combining the cells together to produce the library). The library can be composed of multiple types of cells containing a single or multiple constructs which were integrated into the  
5 genome of a cell treated with radiation, restriction enzymes, and/or DNA breaking agents, applied either together (to the same cells) or separately (applied to individual groups of cells and then combining the cells together to produce the library).

The invention is also directed to methods for making libraries by selecting  
10 various subsets of cells from the same or different transfection experiments. For example, all of the cells expressing nuclear factors (as determined by the presence of nuclear green fluorescent protein in cells transfected with construct 20) can be pooled to create a library of cells with activated nuclear factors. Similarly, cells expressing membrane or secreted proteins can be pooled. Cells can also be  
15 grouped by phenotype, for example, growth factor independent growth, growth factor independent proliferation, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell, muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), gain or loss of cell-cell adhesion, migration, or cellular  
20 activation (e.g., resting versus activated T cells).

The invention is also directed to methods of using libraries of cells to over-express an endogenous gene. The library is screened for the expression of the gene and cells are selected that express the desired gene product. The cell can then be used to purify the gene product for subsequent use. Expression of the cell  
25 can occur by culturing the cell *in vitro* or by allowing the cell to express the gene *in vivo*.

The invention is also directed to methods of using libraries to identify novel gene and gene products.

The invention is also directed to methods for increasing the efficiency of  
30 gene activation by treating the cells with agents that stimulate or effect the

patterns of non-homologous integration. It has been demonstrated that gene expression patterns, chromatin structure, and methylation patterns can differ dramatically from cell type to cell type. Even different cell lines from the same cell type can have significant differences. These differences can impact the patterns of non-homologous integration by affecting both the DNA breakage pattern and the repair process. For example, chromatinized stretches of DNA (characteristics likely associated with inactive genes) may be more resistant to breakage by restriction enzymes and chemical agents, whereas they may be susceptible to breakage by radiation.

Furthermore, inactive genes can be methylated. In this case, restriction enzymes that are blocked by CpG methylation will be unable to cleave methylated sites near the inactive gene, making it more difficult to activate that gene using methylation-sensitive enzymes. These problems can be circumvented by creating activation libraries in several cell lines using a variety of DNA breakage agents. By doing this, a more complete integration pattern can be created and the probability of activating a given gene maximized.

The methods of the invention can include introducing double strand breaks into the DNA of the cell containing the endogenous gene to be over-expressed. These methods introduce double-strand breaks into the genomic DNA in the cell prior to or simultaneously with vector integration. The mechanism of DNA breakage can have a significant effect on the pattern of DNA breaks in the genome. As a result, DNA breaks produced spontaneously or artificially with radiation, restriction enzymes, bleomycin, or other breaking agents, can occur in different locations.

In order to increase integration efficiency and to improve the random distribution of integration sites, cells can be treated with low, intermediate, or high doses of radiation prior to or following transfection. By artificially inducing double strand breaks, the transfected DNA can now integrate into the host cell chromosome as part of the DNA repair process. Normally, creation of double strand breaks to serve as the site of integration is the rate limiting step. Thus, by

increasing chromosome breaks using radiation (or other DNA damaging agents), a larger number of integrants can be obtained in a given transfection. Furthermore, the mechanism of DNA breakage by radiation is different than by spontaneous breakage.

5                Radiation can induce DNA breaks directly when a high energy photon hits the DNA molecule. Alternatively, radiation can activate compounds in the cell which in turn, react with and break the DNA strand. Spontaneous breaks, on the other hand, are thought to occur by the interaction between reactive compounds produced in the cell (such as superoxides and peroxides) and the DNA molecule.  
10              However, DNA in the cell is not present as a naked, deproteinized polymer, but instead is bound to chromatin and present in a condensed state. As a result, some regions are not accessible to agents in the cell that cause double strand breaks. The photons produced by radiation have wave lengths short enough to hit highly condensed regions of DNA, thereby inducing breaks in DNA regions that are  
15              under represented in spontaneous breaks. Thus, radiation is capable of creating different DNA breakage patterns, which in turn, should lead to different integration patterns.

                 As a result, libraries produced using the same activation construct in cells with and without radiation treatment will potentially contain different sets of  
20              activated genes. Finally, radiation treatment increases efficiency of nonhomologous integration by up to 5-10 fold, allowing complete libraries to be created using fewer cells. Thus, radiation treatment increases the efficiency of gene activation and generates new integration and activation patterns in transfected cells. Useful types of radiation include  $\alpha$ ,  $\beta$ ,  $\gamma$ , x-ray, and ultraviolet  
25              radiation. Useful doses of radiation vary for different cell types, but in general, dose ranges resulting in cell viabilities of 0.1% to >99% are useful. For HT1080 cells, this corresponds to radiation doses from a Cs-137 source of approximately 0.1 rads to 1000 rads. Other doses may also be useful as long as the dose either  
                 increases the integration frequency or changes the pattern of integration sites.

In addition to radiation, restriction enzymes can be used to artificially induce chromosome breaks in transfected cells. As with radiation, DNA restriction enzymes can create chromosome breaks which, in turn, serve as integration sites for the transfected DNA. This larger number of DNA breaks increases the overall efficiency of integration of the activation construct. Furthermore, the mechanism of breakage by restriction enzymes differs from that by radiation, the pattern of chromosome breaks is also likely to be different.

Restriction enzymes are relatively large molecules compared to photons and small metabolites capable of damaging DNA. As a result, restriction enzymes will tend to break regions that are less condensed than the genome as a whole. If the gene of interest lies within an accessible region of the genome, then treatment of the cells with a restriction enzyme can increase the probability of integrating the activation construct upstream of the gene of interest. Since restriction enzymes recognize specific sequences, and since a given restriction site may not lie upstream of the gene of interest, a variety of restriction enzymes can be used. It may also be important to use a variety of restriction enzymes since each enzyme has different properties (e.g., size, stability, ability to cleave methylated sites, and optimal reaction conditions) that affect which sites in the host chromosome will be cleaved. Each enzyme, due to the different distribution of cleavable restriction sites, will create a different integration pattern.

Therefore, introduction of restriction enzymes (or plasmids capable of expressing restriction enzymes) before, during, or after introduction of the activation construct will result in the activation of different sets of genes. Finally, restriction enzyme-induced breaks increase the integration efficiency by up to 5-10 fold (Yorifuji *et al.*, *Mut. Res.* 243:121 (1990)), allowing fewer cells to be transfected to produce a complete library. Thus, restriction enzymes can be used to create new integration patterns, allowing activation of genes which failed to be activated in libraries produced by non-homologous recombination at spontaneous breaks or at other artificially induced breaks.

Restriction enzymes can also be used to bias integration of the activation construct to a desired site in the genome. For example, several rare restriction enzymes have been described which cleave eukaryotic DNA every 50-1000 kilobases, on average. If a rare restriction recognition sequence happens to be located upstream of a gene of interest, by introducing the restriction enzyme at the time of transfection along with the activation construct, DNA breaks can be preferentially upstream of the gene of interest. These breaks can then serve as sites for integration of the activation construct. Any enzyme can be that cleaves in an appropriate location in or near the gene of interest and its site is under-represented in the rest of the genome or its site is over-represented near genes (e.g., restriction sites containing CpG). For genes that have not been previously identified, restriction enzymes with 8 bp recognition sites (e.g., *NotI*, *SfiI*, *PmeI*, *SwaI*, *SseI*, *SrfI*, *SgrAI*, *PacI*, *AscI*, *SgfI*, and *Sse8387I*), enzymes recognizing CpG containing sites (e.g., *EagI*, *Bsi-WI*, *MluI*, and *BssHII*) and other rare cutting enzymes can be used.

In this way, "biased" libraries can be created which are enriched for certain types of activated genes. In this respect, restriction enzyme sites containing CpG dinucleotides are particularly useful since these sites are under-represented in the genome at large, but over-represented in the form of CpG islands at the 5' end of many genes, the very location that is useful for gene activation. Enzymes recognizing these sites, therefore, will preferentially cleave at the 5' end of genic sequences.

Restriction enzymes can be introduced into the host cell by several methods. First, restriction enzymes can be introduced into the cell by electroporation (Yorifuji *et al.*, *Mut. Res.* 243:121 (1990); Winegar *et al.*, *Mut. Res.* 225:49 (1989)). In general, the amount of restriction enzyme introduced into the cell is proportional to its concentration in the electroporation media. The pulse conditions must be optimized for each cell line by adjusting the voltage, capacitance, and resistance. Second, the restriction enzyme can be expressed transiently from a plasmid encoding the enzyme under the control of eukaryotic

regulatory elements. The level of enzyme produced can be controlled by using inducible promoters, and varying the strength of induction. In some cases, it may be desirable to limit the amount of restriction enzyme produced (due to its toxicity). In these cases, weak or mutant promoters, splice sites, translation start  
5 codons, and poly A tails can be utilized to lower the amount of restriction enzyme produced. Third, restriction enzymes can be introduced by agents that fuse with or permeabilize the cell membrane. Liposomes and streptolysin O (Pimplikar *et al.*, *J. Cell Biol.* 125:1025 (1994)) are examples of this type of agent. Finally, mechanical perforation (Beckers *et al.*, *Cell* 50:523-534 (1987)) and  
10 microinjection can also be used to introduce nucleases and other proteins into cells. However, any method capable of delivering active enzymes to a living cell is suitable.

DNA breaks induced by bleomycin and other DNA damaging agents can also produce DNA breakage patterns that are different. Thus, any agent or  
15 incubation condition capable of generating double strand breaks in cells is useful for increasing the efficiency and/or altering the sites of non-homologous recombination. Examples of classes of chemical DNA breaking agents include, but are not limited to, peroxides and other free radical generating compounds, alkylating agents, topoisomerase inhibitors, anti-neoplastic drugs, acids,  
20 substituted nucleotides, and enediynes antibiotics.

Specific chemical DNA breaking agents include, but are not limited to, bleomycin, hydrogen peroxide, cumene hydroperoxide, tert-butyl hydroperoxide, hypochlorous acid (reacted with aniline, 1-naphthylamine or 1 naphthol), nitric acid, phosphoric acid, doxorubicin, 9-deoxydoxorubicin,  
25 demethyl-6-deoxyrubicin, 5-iminodaunorubicin, adriamycin, 4'-(9-acridinylamino)methanesulfon-m-anisidide, neocarzinostatin, 8-methoxycaffeine, etoposide, ellipticine, iododeoxyuridine, and bromodeoxyuridine.

It has been shown that DNA repair machinery in the cell can be induced by pre-exposing the cell to low doses of a DNA breaking agent such as radiation  
30 or bleomycin. By pretreating cells with these agents approximately 24 hours prior



to transfection, the cell will be more efficient at repairing DNA breaks and integrating DNA following transfection. In addition, higher doses of radiation or other DNA breaking agents can be used since the LD50 (the dose that results in lethality in 50% of the exposed cells) is higher following pretreatment. This  
5 allows random activation libraries to be created at multiple doses and results in a different distribution of integration sites within the host cell's chromosomes.

### *Screening*

Once an activation library (or libraries) is created, it can be screened using a number of assays. Depending on the characteristics of the protein(s) of interest  
10 (e.g., secreted versus intracellular proteins) and the nature of the activation construct used to create the library, any or all of the assays described below can be utilized. Other assay formats can also be used.

**ELISA.** Activated proteins can be detected using the enzyme-linked immunosorbent assay. If the activated gene product is secreted, culture  
15 supernatants from pools of activation library cells are incubated in wells containing bound antibody specific for the protein of interest. If a cell or group of cells has activated the gene of interest, then the protein will be secreted into the culture media. By screening pools of library clones (the pools can be from 1 to greater than 100,000 library members), pools containing a cell(s) that has  
20 activated the gene of interest can be identified. The cell of interest can then be purified away from the other library members by sib selection, limiting dilution, or other techniques known in the art. In addition to secreted proteins, ELISA can be used to screen for cells expressing intracellular and membrane-bound proteins. In these cases, instead of screening culture supernatants, a small number of cells  
25 is removed from the library pool (each cell is represented at least 100-1000 times in each pool), lysed, clarified, and added to the antibody-coated wells.

**ELISA Spot Assay.** ELISA spot are coated with antibodies specific for the protein of interest. Following coating, the wells are blocked with 1%

BSA/PBS for 1 hour at 37°C. Following blocking, 100,000 to 500,000 cells from the random activation library are applied to each well (representing ~10% of the total pool). In general, one pool is applied to each well. If the frequency of a cell expressing the protein of interest is 1 in 10,000 (i.e., the pool consists of 10,000 individual clones, one of which expresses the protein of interest), then plating 500,000 cells per well will yield 50 specific cells. Cells are incubated in the wells at 37°C for 24 to 48 hours without being moved or disturbed. At the end of the incubation, the cells are removed and the plate is washed 3 times with PBS/0.05% Tween 20 and 3 times with PBS/1%BSA. Secondary antibodies are applied to the wells at the appropriate concentration and incubated for 2 hours at room temperature or 16 hours at 4°C. These antibodies can be biotinylated or labeled directly with horseradish peroxidase (HRP). The secondary antibodies are removed and the plate is washed with PBS/1% BSA. The tertiary antibody or streptavidin labeled with HRP is added and incubated for 1 hour at room temperature.

**FACS assay.** The fluorescence-activated cell sorter can be used to screen the random activation library in a number of ways. If the gene of interest encodes a cell surface protein, then fluorescently-labeled antibodies are incubated with cells from the activation library. If the gene of interest encodes a secreted protein, then cells can be biotinylated and incubated with streptavidin conjugated to an antibody specific to the protein of interest (Manz *et al.*, *Proc. Natl. Acad. Sci. (USA)* 92:1921 (1995)). Following incubation, the cells are placed in a high concentration of gelatin (or other polymer such as agarose or methylcellulose) to limit diffusion of the secreted protein. As protein is secreted by the cell, it is captured by the antibody bound to the cell surface. The presence of the protein of interest is then detected by a second antibody which is fluorescently labeled. For both secreted and membrane bound proteins, the cells can then be sorted according to their fluorescence signal. Fluorescent cells can then be isolated, expanded, and further enriched by FACS, limiting dilution, or other cell purification techniques known in the art.

**Magnetic Bead Separation.** The principle of this technique is similar to FACS. Membrane bound proteins and captured secreted proteins (as described above) are detected by incubating the activation library with an antibody-conjugated magnetic beads that are specific for the protein of interest. If the protein is present on the surface of a cell, the magnetic beads will bind to that cell. Using a magnet, the cells expressing the protein of interest can be purified away from the other cells in the library. The cells are then released from the beads, expanded, analyzed, and further purified if necessary.

**RT-PCR.** A small number of cells (equivalent to at least the number of individual clones in the pool) is harvested and lysed to allow purification of the RNA. Following isolation, the RNA is reversed-transcribed using reverse transcriptase. PCR is then carried out using primers specific for the cDNA of the gene of interest.

Alternatively, primers can be used that span the synthetic exon in the activation construct and the exon of the endogenous gene. This primer will not hybridize to and amplify the endogenously expressed gene of interest. Conversely, if the activation construct has integrated upstream of the gene of interest and activated gene expression, then this primer, in conjunction with a second primer specific for the gene will amplify the activated gene by virtue of the presence of the synthetic exon spliced onto the exon from the endogenous gene. Thus, this method can be used to detect activated genes in cells that normally express the gene of interest at lower than desired levels.

**Phenotypic Section.** In this embodiment, cells can be selected based on a phenotype conferred by the activated gene. Examples of phenotypes that can be selected for include proliferation, growth factor independent growth, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell, muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), gain or loss of cell-cell adhesion, migration, and cellular activation (e.g., resting versus activated T cells). Isolation of activated cells demonstrating a phenotype, such

as those described above, is important because the activation of an endogenous gene by the integrated construct is presumably responsible for the observed cellular phenotype. Thus, the activated gene may be an important therapeutic drug or drug target for treating or inducing the observed phenotype.

5           The sensitivity of each of the above assays can be effectively increased by transiently upregulating gene expression in the library cells. This can be accomplished for NF- $\kappa$ B site-containing promoters (on the activation construct) by adding PMA and tumor necrosis factor- $\alpha$ , *e.g.*, to the library. Separately, or in conjunction with PMA and TNF- $\alpha$ , sodium butyrate can be added to further  
10       enhance gene expression. Addition of these reagents can increase expression of the protein of interest, thereby allowing a lower sensitivity assay to be used to identify the gene activated cell of interest.

          Since large activation libraries are created to maximize activation of many genes, it is advantageous to organize the library clones in pools. Each pool can  
15       consist of 1 to greater than 100,000 individual clones. Thus, in a given pool, many activated proteins are produced, often in dilute concentrations (due to the overall size of the pool and the limited number of cells within the pool that produce a given activated protein). Thus, concentration of the proteins prior to screening effectively increases the ability to detect the activated proteins in the  
20       screening assay. One particularly useful method of concentration is ultrafiltration; however, other methods can also be used. For example, proteins can be concentrated non-specifically, or semi-specifically by adsorption onto ion exchange, hydrophobic, dye, hydroxyapatite, lectin, and other suitable resins under conditions that bind most or all proteins present. The bound proteins can  
25       then be removed in a small volume prior to screening. It is advantageous to grow the cells in serum free media to facilitate the concentration of proteins.

          In another embodiment, a useful sequence that can be included on the activation construct is an epitope tag. The epitope tag can consist of an amino acid sequence that allows affinity purification of the activated protein (*e.g.*, on  
30       immunoaffinity or chelating matrices). Thus, by including an epitope tag on the

activation construct, all of the activated proteins from an activation library can be purified. By purifying the activated proteins away from other cellular and media proteins, screening for novel proteins and enzyme activities can be facilitated. In some instances, it may be desirable to remove the epitope tag following purification of the activated protein. This can be accomplished by including a protease recognition sequence (e.g., Factor IIa or enterokinase cleavage site) downstream from the epitope tag on the activation construct. Incubation of the purified, activated protein(s) with the appropriate protease will release the epitope tag from the proteins(s).

In libraries in which an epitope tag sequence is located on the activation construct, all of the activated proteins can be purified away from all other cellular and media proteins using affinity purification. This not only concentrates the activated proteins, but also purifies them away from other activities that can interfere with the assay used to screen the library.

Once a pool of clones containing cells over-expressing the gene of interest is identified, steps can be taken to isolate the activated cell. Isolation of the activated cell can be accomplished by a variety of methods known in the art. Examples of cell purification methods include limiting dilution, fluorescence activated cell sorting, magnetic bead separation, sib selection, and single colony purification using cloning rings.

In preferred embodiments of the invention, the methods include a process wherein the expression product is purified. In highly preferred embodiments, the cells expressing the endogenous gene product are cultured so as to produce amounts of gene product feasible for commercial application, and especially diagnostic and therapeutic and drug discovery uses.

Any vector used in the methods described herein can include an amplifiable marker. Thereby, amplification of both the vector and the DNA of interest (i.e., containing the over-expressed gene) occurs in the cell, and further enhanced expression of the endogenous gene is obtained. Accordingly, methods can include a step in which the endogenous gene is amplified.

Once the activated cell has been isolated, expression can be further increased by amplifying the locus containing both the gene of interest and the activation construct. This can be accomplished by each of the methods described below, either separately or in combination.

5           Amplifiable markers are genes that can be selected for higher copy number. Examples of amplifiable markers include dihydrofolate reductase, adenosine deaminase, aspartate transcarbamylase, dihydro-orotase, and carbamyl phosphate synthase. For these examples, the elevated copy number of the  
10           amplifiable marker and flanking sequences (including the gene of interest) can be selected for using a drug or toxic metabolite which is acted upon by the amplifiable marker. In general, as the drug or toxic metabolite concentration increases, cells containing fewer copies of the amplifiable marker die, whereas cells containing increased copies of the marker survive and form colonies. These colonies can be isolated, expanded, and analysed for increased levels of  
15           production of the gene of interest.

          Placement of an amplifiable marker on the activation construct results in the juxtaposition of the gene of interest and the amplifiable marker in the activated cell. Selection for activated cells containing increased copy number of the amplifiable marker and gene of interest can be achieved by growing the cells  
20           in the presence of increasing amounts of selective agent (usually a drug or metabolite). For example, amplification of dihydrofolate reductase (DHFR) can be selected using methotrexate.

          As drug-resistant colonies arise at each increasing drug concentration, individual colonies can be selected and characterized for copy number of the  
25           amplifiable marker and gene of interest, and analyzed for expression of the gene of interest. Individual colonies with the highest levels of activated gene expression can be selected for further amplification in higher drug concentrations. At the highest drug concentrations, the clones will express greatly increased amounts of the protein of interest.

When amplifying DHFR, it is convenient to plate approximately  $1 \times 10^7$  cells at several different concentrations of methotrexate. Useful initial concentrations of methotrexate range from approximately 5 nM to 100 nM. However, the optimal concentration of methotrexate must be determined empirically for each cell line and integration site. Following growth in methotrexate containing media, colonies from the highest concentration of methotrexate are picked and analyzed for increased expression of the gene of interest. The clone(s) with the highest concentration of methotrexate are then grown in higher concentrations of methotrexate to select for further amplification of DHFR and the gene of interest. Methotrexate concentrations in the micromolar and millimolar range can be used for clones containing the highest degree of gene amplification.

Placement of a viral origin of replication(s) (e.g., ori P or SV40 in human cells, and polyoma ori in mouse cells) on the activation construct will result in the juxtaposition of the gene of interest and the viral origin of replication in the activated cell. The origin and flanking sequences can then be amplified by introducing the viral replication protein(s) in trans. For example, when ori P (the origin of replication on Epstein-Barr virus) is utilized, EBNA-1 can be expressed transiently or stably. EBNA-1 will initiate replication from the integrated ori P locus. The replication will extend from the origin bi-directionally. As each replication product is created, it too can initiate replication. As a result, many copies of the viral origin and flanking genomic sequences including the gene of interest are created. This higher copy number allows the cells to produce larger amounts of the gene of interest.

At some frequency, the replication product will recombine to form a circular molecule containing flanking genomic sequences, including the gene of interest. Cells that contain circular molecules with the gene of interest can be isolated by single cell cloning and analysis by Hirt extraction and Southern blotting. Once purified, the cell containing the episomal genomic locus at elevated copy number (typically 10-50 copies) can be propagated in culture. To

achieve higher amplification, the episome can be further boosted by including a second origin adjacent to the first in the original construct. For example, T antigen can be used to boost the copy number of ori P/SV40 episomes to a copy number of ~1000 (Heinzel *et al.*, *J. Virol.* 62:3738 (1988)). This substantial increase in copy number can dramatically increase protein expression.

The invention encompasses over-expression of endogenous genes both *in vivo* and *in vitro*. Therefore, the cells could be used *in vitro* to produce desired amounts of a gene product or could be used *in vivo* to provide that gene product in the intact animal.

The invention also encompasses the proteins produced by the methods described herein. The proteins can be produced from either known, or previously unknown genes. Examples of known proteins that can be produced by this method include, but are not limited to, erythropoietin, insulin, growth hormone, glucocerebrosidase, tissue plasminogen activator, granulocyte-colony stimulating factor, granulocyte/macrophage colony stimulating factor, interferon  $\alpha$ , interferon  $\beta$ , interferon  $\gamma$ , interleukin-2, interleukin-6, interleukin-11, interleukin-12, TGF  $\beta$ , blood clotting factor V, blood clotting factor VII, blood clotting factor VIII, blood clotting factor IX, blood clotting factor X, TSH- $\beta$ , bone growth factor 2, bone growth factor-7, tumor necrosis factor, alpha-1 antitrypsin, anti-thrombin III, leukemia inhibitory factor, glucagon, Protein C, protein kinase C, macrophage colony stimulating factor, stem cell factor, follicle stimulating hormone  $\beta$ , urokinase, nerve growth factors, insulin-like growth factors, insulinotropin, parathyroid hormone, lactoferrin, complement inhibitors, platelet derived growth factor, keratinocyte growth factor, neurotrophin-3, thrombopoietin, chorionic gonadotropin, thrombomodulin, alpha glucosidase, epidermal growth factor, FGF, macrophage-colony stimulating factor, interleukin 4, interleukin 10, and cell surface receptors for each of the above growth factors.

Where the protein product from the activated cell is purified, any method of protein purification known in the art may be employed.



## *Examples*

### *Example 1*

#### *Method: Construction of pRIG-1 (FIG. 5)*

Human DHFR was amplified by PCR from cDNA produced from  
5 HT1080 cells by PCR using the primers DHFR-F1  
(5' TCCTTCGAAGCTTGTCATGGTTGGTTCGCTAAACTGCAT 3')  
and DHFR-R1  
(5' AAACCTTAAGATCGATTAATCATTCTTCTCATATACTTCAA), and  
cloned into the T site in pTARGET™ (Promega) to create pTARGET:DHFR.  
10 The RSV promoter was isolated from PREP9 by digestion with *NheI* and *XbaI*  
and inserted into the *NheI* site of pTARGET:DHFR to create pTgT:RSV+DHFR.  
Oligonucleotides JH169 (5' ATCCACCATGGCTACAGGTGAGTACTCG 3')  
and JH170 (5' GATCCGAGTACTCACCTGTAGCCATGGTGGATTAA 3')  
were annealed and inserted into the I-Ppo-I and *NheI* sites of pTgT:RSV+DHFR  
15 to create pTgT:RSV+DHFR+Exl. a 279 bp region corresponding to nucleotides  
230-508 of pBR322 was PCR amplified using primers Tet F1 (5'  
GGCGAGATCTAGCGCTATATGCGTTGATGCAAT 3') and Tet F2 (5'  
GGCCAGATCTGCTACCTTAAGAGAGCCGAAACAAGCGCTCATGAG  
CCCGAA 3'). Amplification products were digested with *BglII* and cloned into  
20 the *BamHI* site of pTgT:RSV+RSV+DHFR+Exl to create pRIG-1.

#### *Transfection – Creation of pRIG-1 Gene Activation Library in HT1080 Cells*

To activate gene expression, a suitable activation construct is selected  
from the group of constructs described above. The selected activation construct  
is then introduced into cells by any transfection method known in the art.  
25 Examples of transfection methods include electroporation, lipofection, calcium

phosphate precipitation, DEAE dextran, and receptor mediated endocytosis. Following introduction into the cells, the DNA is allowed to integrate into the host cell's genome via non-homologous recombination. Integration can occur at spontaneous chromosome breaks or at artificially induced chromosomal breaks.

5           **Method:** Transfection of human cells with pRIG-I.  $2 \times 10^9$  HH1 cells, an HPRT<sup>-</sup> subclone of HT1080 cells, was grown in 150 mm tissue culture plates to 90% confluency. Media was removed from the cells and saved as conditioned media (see below). Cells were removed from the plate by brief incubation with trypsin, added to media/10% fetal bovine serum to neutralize the trypsin, and  
10           pelleted at 1000 rpm in a Jouan centrifuge for 5 minutes. Cells were washed in 1X PBS, counted, and repelleted as above. The cell pellet was resuspended at  $2.5 \times 10^7$  cells/ml final in 1X PBS (Gibco BRL Cat #14200-075). Cells were then exposed to 50 rads of  $\gamma$  irradiation from a Cs<sup>137</sup> source. pRIG-I was linearized with *Bam*HI, purified with phenol/chloroform, precipitated with ethanol, and  
15           resuspended in PBS. Purified and linearized activation construct was added to the cell suspension to produce a final concentration of 40  $\mu$ g/ml. The DNA/irradiated cell mixture was then mixed and 400  $\mu$ l was placed into each 0.4 cm electroporation cuvettes (Biorad). The cuvettes were pulsed at 250 Volts, 600  $\mu$ Farads, 50 Ohms using an electroporation apparatus (Biorad). Following the  
20           electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into  $\alpha$ MEM/10%FBS containing penicillin/streptomycin (Gibco/BRL). The cells were then plated at approximately  $7 \times 10^6$  cells/150 mm plate containing 35 ml  $\alpha$ MEM/10% FBS/penstrep (33% conditioned media/67% fresh media). Following a 24 hour incubation at 37°C, G418 (Gibco/BRL) was  
25           added to each plate to a final concentration of 500  $\mu$ g/ml from a 60 mg/ml stock. After 4 days of selection, the media was replaced with fresh  $\alpha$ MEM/10% FBS/penstrep/500  $\mu$ g/ml G418. The cells were then incubated for another 7-10 days and the culture supernatant assayed for the presence of new protein factors or stored at -80 °C for later analysis. The drug resistant clones can be stored in  
30           liquid nitrogen for later analysis.

## ***Example 2***

### ***Use of ionizing irradiation to increase the frequency and randomness of DNA integration***

**Method:** HH1 cells were harvested at 90% confluency, washed in 1x PBS, and resuspended at a cell concentration of  $7.5 \times 10^6$  cells/ml in 1X PBS. 15  $\mu$ g linearized DNA (pRIG-I) was added to the cells and mixed. 400  $\mu$ l was added to each electroporation cuvette and pulsed at 250 Volts, 600  $\mu$ Farads, 50 Ohms using an electroporation apparatus (Biorad). Following the electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into 2.5 ml  $\alpha$ MEM/10%FBS/1X penstrep. 300  $\mu$ l of cells from each shock were irradiated at 0, 50, 500, and 5000 rads immediately prior to or at either 1 hour or 4 hours post transfection. Immediately following irradiation, the cells were plated onto tissue culture plates in complete medium. At 24 hours post plating, G418 was added to the culture to a final concentration of 500  $\mu$ g/ml. At 7 days post-selection, the culture medium was replaced with fresh complete medium containing 500  $\mu$ g/ml G418. At 10 days post selection, medium was removed from the plate, the colonies were stained with Coomassie Blue/90% methanol/10% acetic acid and colonies with greater than 50 cells were counted.

## ***Example 3***

### ***Use of restriction enzymes to generate random, semi-random, or targeted breaks in the genome***

**Method:** HH1 cells were harvested at 90% confluence, washed in 1x PBS, and resuspended at a cell concentration of  $7.5 \times 10^6$  cells/ml in 1X PBS. To test the efficiency of integration, 15  $\mu$ g linearized DNA (PGK- $\beta$ geo) was added to each 400  $\mu$ l aliquot of cells and mixed. To several aliquots of cells, restriction enzymes *XbaI*, *NotI*, *HindIII*, *Ippol* (10-500 units) were then added to separate

cell/DNA mixture. 400  $\mu$ l was added to each electroporation cuvette and pulsed at 250 Volts, 600  $\mu$ Farads, 50 Ohms using an electroporation apparatus (Biorad). Following the electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into 2.5 ml  $\alpha$ MEM/10%FBS/IX penstrep. 300  $\mu$ l of 2.5 ml total cells from each shock were plated onto tissue culture plates in complete media. At 24 hours post plating, G418 was added to the culture to a final concentration of 600  $\mu$ g/ml. At 7 days post-selection, the media was replaced with fresh complete media containing 600  $\mu$ g/ml G418. At 10 days post selection, media was removed from the plate, the colonies were stained with Coomassie Blue/90% methanol/10% acetic acid and colonies with greater than 50 cells were counted.